

Privacy Preserving PageRank Algorithm By Using Secure Multi-Party Computation

Ferhat Özgür Çatak

TÜBİTAK BİLGEM Cyber Security Institute, Kocaeli, Turkey
ozgur.catak@tubitak.gov.tr

Abstract

In this work, we study the problem of privacy preserving computation on PageRank algorithm. The idea is to enforce the secure multi party computation of the algorithm iteratively using homomorphic encryption based on Paillier scheme. In the proposed PageRank computation, a user encrypt its own graph data using asymmetric encryption method, sends the data set into different parties in a privacy-preserving manner. Each party computes its own encrypted entity, but learns nothing about the data at other parties.

1 Introduction

In this work, we investigate the problem with two or more parties wish to compute the PageRank algorithm [1] on an encrypted graph data set. In this computation model a user doesn't want to disclose private graph data and want to compute compute intensive large scale graph data set.

Graph data sets often contain some private feature information about personal identifying information and also their sensitive relationships [2, 3, 4]. The disclosure risk of the sensitive data arise with the outsourced computation such cloud computing based computation model. Secure multi party computation based privacy preserving techniques are usually adopted to privacy through secret sharing and homomorphic encryption scheme for the original graph data.

In this paper, we investigate how to perform the PageRank computation over the encrypted graph data. The main contributions are

- The work protects the all structure of a graph data without losing the capability to perform PageRank algorithm over it.
- The proposed method is secure in the semi-honest model. Assuming that all parties correctly follow the protocol.

The rest of the paper is organized as follows. Section 2 defines the privacy, secure multi-party computation and semi-honest security model. Section 3 introduces homomorphic encryption, floating point numbers problem and arbitrarily partitioned data. Section 5 presents our method to compute the PageRank values on the encrypted synthetic graph data. Section 6 concludes the paper.

2 Privacy definition

To decide whether a solution achieves the privacy requirements, we need to know the formal definition of privacy. We use a simplified form of the standard definition of security in the static, semi-honest model due to Goldreich [5]. Our other assumption is that all parties follow the protocol which is called semi-honest security model [6, 7]. A formal definition of private two-party computation in the semi-honest model is given below.

Definition 1. (*privacy w.r.t. semi-honest behavior*): Let $f : \{0, 1\}^* \times \{0, 1\}^* \mapsto \{0, 1\}^* \times \{0, 1\}^*$ be probabilistic polynomial-time functionality, where $f_1(x, y)$ (respectively $f_2(x, y)$) denotes first (resp., second) element of $f(x, y)$; and let Π be two-party protocol for computing f . The view of the first (resp., second) party during an execution of Π , denoted $view_1^\Pi$ (resp., $view_2^\Pi$), be $(x, r_1, m_1, \dots, m_t)$ (resp., $(y, r_2, m_1, \dots, m_t)$). r_1 denotes the outcome of the first (resp., r_2 second) party's internal coin tosses, and m_i denotes the i^{th} message it has received.

The output of the first (resp., second) party during an execution of Π on (x, y) , denoted $OUTPUT_1^\Pi(x, y)$ (resp., $OUTPUT_2^\Pi(x, y)$) is implicit in party's view of the execution.

We say that Π privately computes f , if there exists polynomial-time algorithms, denoted S_1 and S_2 , such that

$$\{(S_1(x, f_1(x, y)), f_2(x, y))\}_{x, y \in \{0, 1\}^*} \stackrel{c}{=} \{VIEW_1^\Pi(x, y)\}_{x, y \in \{0, 1\}^*}$$

$$\{(S_2(x, f_1(x, y)), f_2(x, y))\}_{x, y \in \{0, 1\}^*} \stackrel{c}{=} \{VIEW_2^\Pi(x, y)\}_{x, y \in \{0, 1\}^*}$$

In this section formal definition of privacy is given in the graph data set that can model a variety of security protocols and attacks. The proposed model involves a graph $G(V, E)$, a set of k parties P_1, \dots, P_k . Given a matrix D , let $|D|$ denotes the number of rows in D .

We continue with the definitions for security in the semi-honest model. In semi-honest security model, Assuming that a party correctly follow the protocols using its correct input.

3 Preliminaries

In this section, we briefly introduce preliminary knowledge of homomorphic encryption, floating point numbers and arbitrarily partitioned data.

3.1 Homomorphic Encryption

Homomorphic encryption enables operations on plaintexts to be performed on their respective ciphertexts without disclosing the plaintexts when data is divided between two or more servers as it facilitates computations with ciphertexts. A public-key encryption scheme is additively homomorphic if, given two encrypted messages $Enc(a)$ and $Enc(b)$, there exists a public-key operation \oplus such that $Enc(a) \oplus Enc(b)$ is an encryption of $a + b$. Formally, a cryptosystem is additively homomorphic if for any secret key, public key (sk, pk) the plaintext space $\mathcal{P} = \mathbb{Z}_N$ for $x, y \in \mathbb{Z}_N$.

$$\begin{aligned} E_{pk}(x + y \bmod N) &= E_{pk}(x) \cdot E_{pk}(y) \\ E_{pk}(x \cdot y \bmod N) &= E_{pk}(x)^y \end{aligned} \tag{1}$$

3.1.1 Paillier's Encryption Scheme

Paillier cryptosystem [8] is a probabilistic asymmetric algorithm based on the problem to decide whether a number is an n th residue modulo n^2 [9]. The problem of computing n th residue classes is believed to be computationally hard where n is the product of two large primes.

Given a set of possible plaintexts M , a set of key pairs $K = PK \times SK$ where PK is the public key, SK is the secret key. Paillier homomorphic encryption cryptosystem satisfies the following property of any two plaintexts m_1 and m_2 and a constant value a .

$$D_{sk}(E_{pk}(m_1) \times E_{pk}(m_2)) = m_1 + m_2 \tag{2}$$

$$D_{sk}(E_{pk}(m_1)^a) = a \times m_1 \tag{3}$$

Definition 1. *If the adjacency matrix $A(G)$ of a graph G has no loops, then all entries on the main diagonal of $A(G)$ are zeros. Hence $a_{ij} = 0$ whenever $i = j \forall i, j$.*

The adjacency matrix $A(G)$ of a graph G is sparse and often binary. In order to prevent guessing elements of input data set, Paillier cryptosystem has probabilistic encryption that does not encrypt two equal plain text with the same encryption key into the same ciphertext.

3.2 Floating Point Numbers

Although the proposed protocol manipulate integers, the PageRank algorithm is typically applied to continuous data. However, in the case of real number inputs to the protocol, we need to map data vectors into the discrete domain [10].

Let $ConvertInteger : \mathbb{R}^m \rightarrow \mathbb{Z}^m$ be the corresponding function that multiplies its argument by an exponent $(K : 2^K)$ then rounds them to the nearest

integer value and thus support finite precision. Equation 4 shows the conversion function.

$$\hat{\mathbf{x}} \leftarrow \text{ConvertInteger}(\mathbf{x}) \text{ where } \mathbf{x} \in \mathbb{R}^m, \hat{\mathbf{x}} \in \mathbb{Z}^m \quad (4)$$

3.3 Arbitrarily Partitioned Data

In this work, arbitrary partitioned data between multi-parties (K), $K > 2$ is considered. In the arbitrary partitioned data scheme, there is no specific order of how the data is divided between multiple parties. Specifically, if we have a data set $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, consisting of n row, and each rows in X contains m numeric attributes $\mathbf{x}_i = \{x_i^1 \dots x_i^m\}$. X_i^j is the subset of data set owned by party P_j , then we have $X_i^1 \cup X_i^2 \dots \cup X_i^K = X_i$ and $X_i^1 \cap X_i^2 \dots \cap X_i^K = \emptyset$. In each row (X_i), party P_k has a number of attributes t_i^k , where $\sum_{p=1}^K t_i^k = m$ and each party's attribute size does not have to be equal. If a party has the same attributes in each row, then the arbitrary partition becomes a vertical partition.

4 Privacy-Preserving PageRank

We now formally define the problem. Let k be the number of parties, each having different attributes for the same set of entities (i.e. adjacency matrix rows). Parties P_1, \dots, P_k have as their respective private input sets (i.e., adjacency matrix) S_1, \dots, S_k drawn from some finite universe U . The parties try to measure the importance of nodes with their joint data using the PageRank algorithm in a privacy-preserving manner, i.e., without leaking their elements of S .

4.1 Notation

A directed graph $\mathcal{G}(\mathcal{V}, \mathcal{E}, m)$ representing a web sites set is a set on n nodes connected by a set of m links, where \mathcal{V} denotes the set of nodes and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of links and $m : \mathcal{V} \rightarrow (0, \infty)$ is a weight. For a node $u \in \mathcal{V}$, \mathcal{N} denotes the number of hyperlinks contained in page (node) u .

Given $s_1, s_2 \in \mathcal{V}$, one can say that x and y are neighbors (i.e. hyperlinked from x to y) if $\mathcal{E}(x, y) > 0$ [11].

4.2 Security model

In this paper, the aim is to enable multiple parties (or servers) to jointly conduct the PageRank computation without revealing their private data. Our main assumption is that the input data set is divided between two or more parties, that are willing to compute the PageRank of \mathcal{G} if nothing beyond the expected end results are revealed [12]. Our other assumption is that all parties follow the protocol which is called semi-honest security model [6, 7].

4.3 Graph Representation

In order to protect the unauthorized access to the sensitive graph information, one needs to represent and encrypt the graph in a proper way. In order to achieve these requirements, adjacency matrix representation is used. let k be the number of parties, each has the same attributes for the different entities. The distributed parties try to calculate PageRank algorithm by using their joint data (*i.e.*, adjacency matrix).

The PageRank algorithm needs adjacency matrix \mathcal{A} of graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$. The proposed protocols directly follow the standard the PageRank algorithm. The approximations to the page rank values are iteratively computed until the changes in values in one iteration is below a threshold or max iteration number is reached. At each iteration, each party multiplies its own encrypted adjacency matrix $\llbracket \mathcal{A} \rrbracket_i$ with *PageRank* and *out degree* vectors as shown in Equation 5.

$$\llbracket \mathcal{T} \rrbracket_i = \begin{bmatrix} \llbracket a \rrbracket_{1,1} & \cdots & \llbracket a \rrbracket_{1,j} \\ \vdots & \ddots & \vdots \\ \llbracket a \rrbracket_{m,1} & \cdots & \llbracket a \rrbracket_{m,j} \end{bmatrix} \begin{bmatrix} Pr_1 \\ \vdots \\ Pr_m \end{bmatrix} \begin{bmatrix} \frac{1}{d_1} \\ \vdots \\ \frac{1}{d_m} \end{bmatrix} \quad (5)$$

where $\llbracket a \rrbracket$ is the encrypted elements of $\llbracket \mathcal{A} \rrbracket$, Pr is the PageRank values of graph \mathcal{G} and d is the out degree values. Once the $\llbracket \mathcal{T} \rrbracket_i$ are computed at each party, then another trusted party merges the matrices to find a global encrypted $\llbracket \mathcal{T} \rrbracket$ matrix as shown in Equation 6.

$$\llbracket \mathcal{T} \rrbracket = \bigcup_{i=1}^k \llbracket \mathcal{T} \rrbracket_i \quad (6)$$

The trusted party decrypts $\llbracket \mathcal{T} \rrbracket$ using private key key_{priv} , calculates the PageRank values as shown in Equation 7.

$$Pr(\mathcal{V}_i) = \frac{1-d}{\mathcal{N}} + d \sum_{p_j \in M(p_i)} \frac{Pr(\mathcal{V}_j)}{L(\mathcal{V}_j)} \quad (7)$$

At each iteration, the PageRank is calculated, *i.e.*, the proposed protocols securely find the PageRank vector for graph \mathcal{G} . Once the initial PageRank values are known, the next Pagerank values can be computed locally. Protocol 1 use Protocol 2 (*secureIntMatrix*) to compute secure intermediate matrix.

The proposed algorithm is shown in Protocol 1 - 2.

5 Experiments

We evaluated the proposed algorithm on a synthetic graph data set with 20 nodes. We implemented the protocols in Python 2.7 with the Paillier¹ library. The experiments were performed on a computer with a 2.6 GHz Intel Core

¹<https://github.com/mikeivanov/paillier>

Input: Adjacency matrix: $\mathcal{A} \in \mathbb{R}^{m \times m}$ party size: k , Crypto key length: l , damping factor d ,
Result: Encrypted sub data sets for each party P_s

```

begin
   $\mathcal{D} \leftarrow \text{normalize}(\mathcal{D});$ 
   $Key_{pub}, Key_{priv} \leftarrow KeyGen(l)$   $\triangleright$  Generate public/private keys ;
  for  $i \in k$  do
     $// \text{ create sub adjacency matrix } \mathcal{A}_i \text{ with random feature index for party } P_i ;$ 
     $[\![\mathcal{A}_i]\!] \leftarrow \text{encrypt}(\mathcal{A}^i, Key_{pub});$ 
     $\text{sendToParty}([\![\mathcal{A}]\!]_i, Key_{pub});$ 
  end
  repeat
    foreach  $P_s \in P$  do
       $\text{call } \text{secureIntMatrix}([\![\mathcal{A}]\!]_s) \text{ \{Protocol 2\}}$   $\triangleright$  for each party;
    end
     $comp\_mat \leftarrow \bigcup_{i=1}^k \text{decrypt}([\![\mathcal{T}]\!]_i, key_{priv})$   $\triangleright$  Decrypt, merge;
    for  $i = 1 \dots m$  do
       $PageRank[i] \leftarrow \frac{1-d}{N} + d \sum_{j=1}^m comp\_mat[j, i]$ 
    end
  until Pagerank converges;
end

```

Protocol 1: Adjacency matrix split and encryption

i5 processor and 4 GB main memory running Mac OSX. The execution time is measured by used processor time in seconds with different node size and different encryption key length. Experimental results are shown in Table 1.

6 Conclusion and Future Works

We introduce a new homomorphic encryption based privacy-preserving PageRank algorithm with secure multi-party computation approach. We showed the effects of different encryption key length on results with graphics. In particular, different party size and key length could also be investigated for both in computation time and computation error because of scaling problem of the encryption method.

References

- [1] L. Page, S. Brin, R. Motwani, and T. Winograd, “The pagerank citation ranking: Bringing order to the web,” in *Proceedings of the 7th International World Wide Web Conference*, (Brisbane, Australia), pp. 161–172, 1998.

Input: Encrypted sub adjacency matrix: $\llbracket \mathcal{A} \rrbracket_s \in \mathbb{R}^{m \times n}$, Public crypto key: Key_{pub} , scaling factor c
Result: Intermediate results of *computation_matrix_s*
begin
 $\llbracket comp_mat \rrbracket_s = \text{zeros}(m, n)$;
 for $i = 1 \dots n$ **do**
 $\text{col} \leftarrow \llbracket \mathcal{A} \rrbracket_s[:, i]$ \triangleright Get i^{th} column vector of $\llbracket \mathcal{A} \rrbracket_s$;
 for $j = 1 \dots m$ **do**
 $\llbracket comp_mat \rrbracket_s[i, j] = \llbracket \mathcal{A} \rrbracket_s[i, j] \otimes \left(\frac{PageRank[i]}{OutDegree[i]} \cdot 10^c \right)$ \triangleright
 Encrypted multiplication ;
 end
 end
 return $\llbracket comp_mat \rrbracket_s$
end

Protocol 2: secureIntMatrix

Table 1: Experimental results

Party Size	128 bit	256 bit	512 bit	1024 bit
3	56.45	107.13	678.29	3711,48
5	57.14	110.14	685.70	3506,84
7	56.94	106.70	703.70	3416,14
10	53.65	108.72	962.98	3427,22

- [2] J. Brickell and V. Shmatikov, “Privacy-preserving graph algorithms in the semi-honest model,” in *Advances in Cryptology-ASIACRYPT 2005*, pp. 236–252, Springer, 2005.
- [3] E. Zheleva and L. Getoor, “Preserving the privacy of sensitive relationships in graph data,” in *Privacy, security, and trust in KDD*, pp. 153–171, Springer, 2008.
- [4] B. Zhou, J. Pei, and W. Luk, “A brief survey on anonymization techniques for privacy preserving publishing of social network data,” *ACM Sigkdd Explorations Newsletter*, vol. 10, no. 2, pp. 12–22, 2008.
- [5] O. Goldreich, S. Micali, and A. Wigderson, “How to play any mental game,” in *Proceedings of the nineteenth annual ACM symposium on Theory of computing*, pp. 218–229, ACM, 1987.
- [6] Y. Lindell and B. Pinkas, “Privacy preserving data mining,” in *Advances in CryptologyCRYPTO 2000*, pp. 36–54, Springer, 2000.

- [7] Z. Yang, S. Zhong, and R. N. Wright, “Privacy-preserving classification of customer data without loss of accuracy,” in *SDM*, pp. 92–102, SIAM, 2005.
- [8] P. Paillier, “Public-key cryptosystems based on composite degree residuosity classes,” in *Advances in cryptologyEUROCRYPT99*, pp. 223–238, Springer, 1999.
- [9] C. Orlandi, A. Piva, and M. Barni, “Oblivious neural network computing via homomorphic encryption,” *EURASIP Journal on Information Security*, vol. 2007, p. 18, 2007.
- [10] L. Kamm and J. Willemson, “Secure floating point arithmetic and private satellite collision analysis,” *International Journal of Information Security*, vol. 14, no. 6, pp. 531–548, 2015.
- [11] H. Baloudi, S. Golenia, and A. Jeribi, “The adjacency matrix and the discrete laplacian acting on forms,” *arXiv preprint arXiv:1505.06109*, 2015.
- [12] A. Bansal, T. Chen, and S. Zhong, “Privacy preserving back-propagation neural network learning over arbitrarily partitioned data,” *Neural Computing and Applications*, vol. 20, no. 1, pp. 143–150, 2011.